



Securing Artificial Intelligence (SAI); Mitigation Strategy Report (standards.iteh.ai)

[ETSI GR SAI 005 V1.1.1 \(2021-03\)](https://standards.iteh.ai/catalog/standards/sist/b68ecdff-bbec-48c0-a490-35a21b959449/etsi-gr-sai-005-v1-1-1-2021-03)
<https://standards.iteh.ai/catalog/standards/sist/b68ecdff-bbec-48c0-a490-35a21b959449/etsi-gr-sai-005-v1-1-1-2021-03>

Disclaimer

The present document has been produced and approved by the Secure AI (SAI) ETSI Industry Specification Group (ISG) and represents the views of those members who participated in this ISG. It does not necessarily represent the views of the entire ETSI membership.

Reference

DGR/SAI-005

Keywords

artificial intelligence, security

ETSI

650 Route des Lucioles
F-06921 Sophia Antipolis Cedex - FRANCE

Tel.: +33 4 92 94 42 00 Fax: +33 4 93 65 47 16

Siret N° 348 623 562 00017 - NAF 742 C
Association à but non lucratif enregistrée à la
Sous-Préfecture de Grasse (06) N° 7803/88

iTeh STANDARD PREVIEW
(standards.iteh.ai)

Important notice

<https://standards.iteh.ai/catalog/standards/sist/b68ecdff-bbec-48c0-a490-35a21b759449/etsi-gr-sai-005-v1-1-1-2021-03>
The present document can be downloaded from:
<http://www.etsi.org/standards-search>

The present document may be made available in electronic versions and/or in print. The content of any electronic and/or print versions of the present document shall not be modified without the prior written authorization of ETSI. In case of any existing or perceived difference in contents between such versions and/or in print, the prevailing version of an ETSI deliverable is the one made publicly available in PDF format at www.etsi.org/deliver.

Users of the present document should be aware that the document may be subject to revision or change of status. Information on the current status of this and other ETSI documents is available at <https://portal.etsi.org/TB/ETSIDeliverableStatus.aspx>

If you find errors in the present document, please send your comment to one of the following services:
<https://portal.etsi.org/People/CommitteeSupportStaff.aspx>

Copyright Notification

No part may be reproduced or utilized in any form or by any means, electronic or mechanical, including photocopying and microfilm except as authorized by written permission of ETSI.

The content of the PDF version shall not be modified without the written authorization of ETSI.

The copyright and the foregoing restriction extend to reproduction in all media.

© ETSI 2021.

All rights reserved.

DECT™, **PLUGTESTS™**, **UMTS™** and the ETSI logo are trademarks of ETSI registered for the benefit of its Members.

3GPP™ and **LTE™** are trademarks of ETSI registered for the benefit of its Members and of the 3GPP Organizational Partners.

oneM2M™ logo is a trademark of ETSI registered for the benefit of its Members and of the oneM2M Partners.

GSM® and the GSM logo are trademarks registered and owned by the GSM Association.

Contents

Intellectual Property Rights	4
Foreword.....	4
Modal verbs terminology.....	4
1 Scope	5
2 References	5
2.1 Normative references	5
2.2 Informative references.....	5
3 Definition of terms, symbols and abbreviations.....	11
3.1 Terms.....	11
3.2 Symbols.....	11
3.3 Abbreviations	12
4 Overview	12
4.1 Machine learning models workflow	12
4.2 Mitigation strategy framework.....	13
5 Mitigations against training attacks.....	14
5.1 Introduction	14
5.2 Mitigating poisoning attacks	15
5.2.1 Overview	15
5.2.2 Model enhancement mitigations against poisoning attacks	15
5.2.3 Model-agnostic mitigations against poisoning attacks.....	16
5.3 Mitigating backdoor attacks	16
5.3.1 Overview	16
5.3.2 Model enhancement mitigations against backdoor attacks	17
5.3.3 Model-agnostic mitigations against backdoor attacks.....	18
6 Mitigations against inference attacks.....	19
6.1 Introduction	19
6.2 Mitigating evasion attacks.....	20
6.2.1 Overview	20
6.2.2 Model enhancement mitigations against evasion attacks.....	20
6.2.3 Model-agnostic mitigations against evasion attacks	23
6.3 Mitigating model stealing.....	24
6.3.1 Overview	24
6.3.2 Model enhancement mitigations against model stealing.....	25
6.3.3 Model-agnostic mitigations against model stealing	26
6.4 Mitigating data extraction	27
6.4.1 Overview	27
6.4.2 Model enhancement mitigations against data extraction	27
6.4.3 Model-agnostic mitigations against data extraction.....	28
7 Conclusion.....	29
Annex A: Change History	30
History	31

Intellectual Property Rights

Essential patents

IPRs essential or potentially essential to normative deliverables may have been declared to ETSI. The information pertaining to these essential IPRs, if any, is publicly available for **ETSI members and non-members**, and can be found in ETSI SR 000 314: *"Intellectual Property Rights (IPRs); Essential, or potentially Essential, IPRs notified to ETSI in respect of ETSI standards"*, which is available from the ETSI Secretariat. Latest updates are available on the ETSI Web server (<https://ipr.etsi.org/>).

Pursuant to the ETSI IPR Policy, no investigation, including IPR searches, has been carried out by ETSI. No guarantee can be given as to the existence of other IPRs not referenced in ETSI SR 000 314 (or the updates on the ETSI Web server) which are, or may be, or may become, essential to the present document.

Trademarks

The present document may include trademarks and/or tradenames which are asserted and/or registered by their owners. ETSI claims no ownership of these except for any which are indicated as being the property of ETSI, and conveys no right to use or reproduce any trademark and/or tradename. Mention of those trademarks in the present document does not constitute an endorsement by ETSI of products, services or organizations associated with those trademarks.

Foreword

This Group Report (GR) has been produced by ETSI Industry Specification Group (ISG) Secure AI (SAI).

Modal verbs terminology

In the present document "should", "should not", "may", "need not", "will", "will not", "can" and "cannot" are to be interpreted as described in clause 3.2 of the [ETSI Drafting Rules](https://standards.ietf.org/rfc/rfc2119/) (Verbal forms for the expression of provisions).

"**must**" and "**must not**" are **NOT** allowed in ETSI deliverables except when used in direct citation.

1 Scope

The present document summarizes and analyses existing and potential mitigation against threats for AI-based systems as discussed in ETSI GR SAI 004 [i.1]. The goal is to have a technical survey for mitigating against threats introduced by adopting AI into systems. The technical survey shed light on available methods of securing AI-based systems by mitigating against known or potential security threats. It also addresses security capabilities, challenges, and limitations when adopting mitigation for AI-based systems in certain potential use cases.

2 References

2.1 Normative references

Normative references are not applicable in the present document.

2.2 Informative references

References are either specific (identified by date of publication and/or edition number or version number) or non-specific. For specific references, only the cited version applies. For non-specific references, the latest version of the referenced document (including any amendments) applies.

NOTE: While any hyperlinks included in this clause were valid at the time of publication, ETSI cannot guarantee their long term validity.

The following referenced documents are not necessary for the application of the present document but they assist the user with regard to a particular subject area.

[i.1] ETSI GR SAI 004: "Securing Artificial Intelligence (SAI); Problem Statement".

[i.2] Doyen Sahoo, Quang Pham, Jing Lu, Steven C. H. Hoi: "Online Deep Learning: Learning Deep Neural Networks on the Fly", International Joint Conferences on Artificial Intelligence Organization, 2018.

NOTE: Available at <https://doi.org/10.24963/ijcai.2018/369>.

[i.3] Battista Biggio and Fabio Roli: "Wild patterns: Ten years after the rise of adversarial machine learning", Pattern Recognition, 2018.

[i.4] Qiang Liu, Pan Li, Wentao Zhao, Wei Cai, Shui Yu and Victor C. M. Leung: "A Survey on Security Threats and Defensive Techniques of Machine Learning: A Data Driving View". IEEE Access 2018.

NOTE: Available at <https://doi.org/10.1109/ACCESS.2018.2805680>.

[i.5] Nicolas Papernot, Patrick D. McDaniel, Arunesh Sinha and Michael P. Wellman: "SoK: Security and Privacy in Machine Learning". IEEE European Symposium on Security and Privacy (EuroS&P) 2018.

[i.6] Han Xu, Yao Ma, Haochen Liu, Debayan Deb, Hui Liu, Jiliang Tang and Anil K. Jain: "Adversarial Attacks and Defenses in Images, Graphs and Text: A Review". International Journal of Automation and Computing volume 17, pages151-178(2020).

NOTE: Available at <https://doi.org/10.1007/s11633-019-1211-x>.

[i.7] Anirban Chakraborty, Manar Alam, Vishal Dey, Anupam Chattopadhyay, Debdeep Mukhopadhyay: "Adversarial Attacks and Defences: A Survey", arXiv preprint arXiv:1810.00069v1.

NOTE: Available at <https://arxiv.org/abs/1810.00069>.

- [i.8] Yingzhe He, Guozhu Meng, Kai Chen, Xingbo Hu, Jinwen He: "Towards Privacy and Security of Deep Learning Systems: A Survey". IEEE Transactions on Software Engineering 2020.
- [i.9] NIST IR 8269-(Draft): "A Taxonomy and Terminology of Adversarial Machine Learning".
- NOTE: Available at <https://doi.org/10.6028/NIST.IR.8269-draft>.
- [i.10] Christian Berghoff¹, Matthias Neu¹ and Arndt von Twickel: "Vulnerabilities of Connectionist AI Applications: Evaluation and Defence", Frontiers in Big Data volume 3, 2020.
- NOTE: Available at <https://doi.org/10.3389/fdata.2020.00023>.
- [i.11] Blaine Nelson, Marco Barreno, Fuching Jack Chi, Anthony D. Joseph, Benjamin I. P. Rubinstein, Udam Saini, Charles A. Sutton, J. Doug Tygar, Kai Xia: "Exploiting Machine Learning to Subvert Your Spam Filter", Usenix Workshop on Large-Scale Exploits and Emergent Threats (LEET) 2008.
- [i.12] Matthew Jagielski, Alina Oprea, Battista Biggio, Chang Liu, Cristina Nita-Rotaru, Bo Li: "Manipulating Machine Learning: Poisoning Attacks and Countermeasures for Regression Learning", IEEE Symposium on Security and Privacy 2018: 19-35.
- [i.13] Nathalie Baracaldo, Bryant Chen, Heiko Ludwig, Jaehoon Amir Safavi: "Mitigating Poisoning Attacks on Machine Learning Models: A Data Provenance Based Approach", AISec@CCS 2017: 103-110.
- NOTE: Available at <https://doi.org/10.1145/3128572.3140450>.
- [i.14] Sanghyun Hong, Varun Chandrasekaran, Yigitcan Kaya, Tudor Dumitras, Nicolas Papernot: "On the Effectiveness of Mitigating Data Poisoning Attacks with Gradient Shaping", arXiv: 2002.11497v2.
- NOTE: Available at <https://arxiv.org/abs/2002.11497v2>.
- [i.15] Nitika Khurana, Sudip Mittal, Aritran Dey, Anupam Joshi: "Preventing Poisoning Attacks On AI Based Threat Intelligence Systems", IEEE International Workshop on Machine Learning for Signal Processing (MLSP) 2019: 1-6.
- NOTE: Available at <https://doi.org/10.1109/MLSP.2019.8918803>.
- [i.16] Battista Biggio, Igino Corona, Giorgio Fumera, Giorgio Giacinto, Fabio Roli: "Bagging Classifiers for Fighting Poisoning Attacks in Adversarial Classification Tasks", International Workshop on Multiple Classifier Systems (MCS) 2011: 350-359.
- NOTE: Available at https://doi.org/10.1007/978-3-642-21557-5_37.
- [i.17] Yao Cheng, Cheng-Kang Chu, Hsiao-Ying Lin, Marius Lombard-Platet, David Naccache: "Keyed Non-parametric Hypothesis Tests", International Conference on Network and System Security (NSS) 2019: 632-645.
- NOTE: Available at https://doi.org/10.1007/978-3-030-36938-5_39.
- [i.18] Tran, Brandon, Jerry Li, and Aleksander Madry: "Spectral signatures in backdoor attacks", In Advances in Neural Information Processing Systems, pp. 8000-8010. 2018.
- [i.19] Chen, Bryant, Wilka Carvalho, Nathalie Baracaldo, Heiko Ludwig, Benjamin Edwards, Taesung Lee, Ian Molloy and Biplav Srivastava: "Detecting backdoor attacks on deep neural networks by activation clustering", Artificial Intelligence Safety Workshop @ AAAI, 2019.
- [i.20] Yuntao Liu, Yang Xie, Ankur Srivastava: "Neural Trojans", 2017 IEEE International Conference on Computer Design (ICCD), Boston, MA, 2017, pp. 45-48, doi: 10.1109/ICCD.2017.16.
- NOTE: Available at <https://doi.org/10.1109/ICCD.2017.16>.

- [i.21] Bingyin Zhao and Yangjie Lao: "Resilience of Pruned Neural Network against poisoning attack", International Conference on Malicious and Unwanted Software (MALWARE) 2018, page 78-83.
- NOTE: <https://doi.org/10.1109/MALWARE.2018.8659362>.
- [i.22] Liu, Kang, Brendan Dolan-Gavitt and Siddharth Garg: "Fine-pruning: Defending against backdooring attacks on deep neural networks", In International Symposium on Research in Attacks, Intrusions and Defenses, pp. 273-294. Springer, Cham, 2018.
- NOTE: Available at https://doi.org/10.1007/978-3-030-00470-5_13.
- [i.23] Wang, Bolun, Yuanshun Yao, Shawn Shan, Huiying Li, Bimal Viswanath, Haitao Zheng and Ben Y. Zhao: "Neural cleanse: Identifying and mitigating backdoor attacks in neural networks", In 2019 IEEE Symposium on Security and Privacy (SP), pp. 707-723.
- [i.24] Wenbo Guo, Lun Wang, Xinyu Xing, Min Du, and Dawn Song. 2019: "Tabor: A highly accurate approach to inspecting and restoring trojan backdoors in ai systems", arXiv preprint arXiv:1908.01763v2 (2019).
- NOTE: Available at <https://arxiv.org/abs/1908.01763v2>.
- [i.25] Yansong Gao, Chang Xu, Derui Wang, Shiping Chen, Damith C. Ranasinghe, Surya Nepal: "STRIP: A Defence Against Trojan Attacks on Deep Neural Networks", 2019 Annual Computer Security Applications Conference (ACSAC '19).
- [i.26] Chou Edward, Florian Tramèr, Giancarlo Pellegrino: "sentiNet: Detecting Localized Universal Attack Against Deep Learning Systems", The 3rd Deep Learning and Security Workshop (2020).
- [i.27] Sakshi Udeshi, Shanshan Peng, Gerald Woo, Lionell Loh, Louth Rawshan and Sudipta Chattopadhyay. 2019: "Model Agnostic Defence against Backdoor Attacks in Machine Learning". arXiv preprint arXiv:1908.02203v2 (2019).
- NOTE: Available at <https://arxiv.org/abs/1908.02203v2>.
- [i.28] Bao Gia Doan, Ehsan Abbasnejad, and Damith Ranasinghe: "Februus: Input Purification Defense Against Trojan Attacks on Deep Neural Network Systems", The 36th Annual Computer Security Applications Conference (ACSAC 2020).
- [i.29] Ramprasaath R. Selvaraju, Michael Cogswell, Abhishek Das, Ramakrishna Vedantam, Devi Parikh, Dhruv Batra: "Grad-CAM: Visual Explanations from Deep Networks via Gradient-based Localization", International Conference on Computer Vision (ICCV'17).
- [i.30] Xiaojun Xu, Qi Wang, Huichen Li, Nikita Borisov, Carl A Gunter and Bo Li: "Detecting AI Trojans Using Meta Neural Analysis", IEEE S&P 2021.
- [i.31] Huili Chen, Cheng Fu, Jishen Zhao, Farinaz Koushanfar: "DeepInspect: A Black-box Trojan Detection and Mitigation Framework for Deep Neural Networks", IJCAI 2019: 4658-4664.
- [i.32] Soheil Kolouri, Aniruddha Saha, Hamed Pirsiavash, Heiko Hoffmann, Universal Litmus Patterns: "Revealing Backdoor Attacks in CNNs", CVPR 2020.
- [i.33] Ambra Demontis, Marco Melis, Maura Pintor, Matthew Jagielski, Battista Biggio, Alina Oprea, Cristina Nita-Rotaru, Fabio Roli: "Why Do Adversarial Attacks Transfer? Explaining Transferability of Evasion and Poisoning Attacks", USENIX Security Symposium 2019.
- [i.34] Yinpeng Dong, Tianyu Pang, Hang Su, Jun Zhu: "Evading Defenses to Transferable Adversarial Examples by Translation-Invariant Attacks", CVPR 2019.
- [i.35] Carlini and Wagner: "Adversarial Examples Are Not Easily Detected: Bypassing Ten Detection Methods", the 10th ACM Workshop on Artificial Intelligence and Security 2017.
- [i.36] Anish Athalye, Nicholas Carlini, David Wagner: "Obfuscated Gradients Give a False Sense of Security: Circumventing Defenses to Adversarial Examples", ICML 2018.
- [i.37] Florian Tramèr, Nicholas Carlini, Wieland Brendel, Aleksander Madry: "On Adaptive Attacks to Adversarial Example Defenses", NIPS 2020.

- [i.38] Florian Tramèr, Jens Behrmann, Nicholas Carlini, Nicolas Papernot, Jörn-Henrik Jacobsen: "Fundamental Tradeoffs between Invariance and Sensitivity to Adversarial Perturbations", ICML 2020.
- [i.39] Gintare Karolina Dziugaite, Zoubin Ghahramani, Daniel M. Roy: "A study of the effect of JPG compression on adversarial images", International Society for Bayesian Analysis (ISBA 2016) World Meeting.
- [i.40] [Uri Shaham](#), [James Garritano](#), [Yutaro Yamada](#), [Ethan Weinberger](#), [Alex Cloninger](#), [Xiuyuan Cheng](#), [Kelly Stanton](#), [Yuval Kluger](#): "Defending against adversarial images using basis functions transformation", ArXiv 2018.

NOTE: Available at <https://arxiv.org/abs/1803.10840v3>.

- [i.41] Richard Shin, Dawn Song: "JPEG-resistant adversarial images", in Proc. Mach. Learn. Comput. Secur. Workshop, 2017.
- [i.42] Nilaksh Das, Madhuri Shanbhogue, Shang-Tse Chen, Fred Hohman, Siwei Li, Li Chen, Michael E. Kounavis, Duen Horng Chau: "Shield: Fast, Practical Defense and Vaccination for Deep Learning using JPEG Compression", KDD 2018.
- [i.43] Chuan Guo, Mayank Rana, Moustapha Cissé, Laurens van der Maaten: "Countering Adversarial Images using Input Transformations", ICLR 2018.
- [i.44] [Hossein Hosseini](#), [Yize Chen](#), [Sreeram Kannan](#), [Baosen Zhang](#), [Radha Poovendran](#): "Blocking Transferability of Adversarial Examples in Black-Box Learning Systems", ArXiv 2017.

NOTE: Available at <https://arxiv.org/abs/1703.04318>.

- [i.45] Ian J. Goodfellow, Jonathon Shlens, Christian Szegedy: "Explaining and Harnessing Adversarial Examples", ICLR 2015.
- [i.46] Aleksander Madry, Aleksandar Makelov, Ludwig Schmidt, Dimitris Tsipras, Adrian Vladu: "Towards Deep Learning Models Resistant to Adversarial Attacks", ICLR 2018.
- [i.47] Madry et al. (2019), <https://standards.iteh.ai/catalog/standards/sist/b68ecdff-bbec-48c0-a490-53a21b959447/etsi-gr-sai-005-v1-1-1-2021-03>, RobustML.

NOTE: Available at <https://www.robust-ml.org/>.

- [i.48] Moustapha Cisse, Piotr Bojanowski, Edouard Grave, Yann Dauphin, Nicolas Usunier: "Parseval Networks: Improving robustness to adversarial examples", ICML 2017.
- [i.49] Ross and Doshi-Velez: "Improving the adversarial robustness and interpretability of deep neural networks by regularizing their input gradients", AAAI 2018.
- [i.50] Cem Anil, James Lucas, Roger Grosse: "Sorting out Lipschitz function approximation", ICML 2019.
- [i.51] Jeremy Cohen, Elan Rosenfeld, Zico Kolter: "Certified Adversarial Robustness via Randomized Smoothing", ICML 2019.
- [i.52] Nicolas Papernot, Patrick McDaniel, Xi Wu, Somesh Jha and Ananthram Swami: "Distillation as a defense to adversarial perturbations against deep neural networks", 2016 IEEE Symposium S&P.
- [i.53] Nicolas Carlini, David Wagner: "Towards evaluating the robustness of neural networks", 2017 IEEE symposium on Security and Privacy (S&P).
- [i.54] Akhilan Boopathy, Tsui-Wei Weng, Pin-Yu Chen, Sijia Liu, Luca Daniel: "CNN-Cert: An Efficient Framework for Certifying Robustness of Convolutional Neural Networks", AAAI 2019.
- [i.55] Ching-Yun Ko, Zhaoyang Lyu, Tsui-Wei Weng, Luca Daniel, Ngai Wong, Dahua Lin: "POPQORN: Quantifying Robustness of Recurrent Neural Networks", ICML 2019.
- [i.56] [Aleksandar Bojchevski](#), [Stephan Günnemann](#): "Certifiable Robustness to Graph Perturbations", NIPS 2019.

- [i.57] G. Katz, C. Barrett, D. L. Dill, K. Julian and M. J. Kochenderfer: "Reluplex: An efficient SMT solver for verifying deep neural networks", in International Conference on Computer Aided Verification. Springer, 2017, pp. 97-117.
- [i.58] X. Huang, M. Kwiatkowska, S. Wang and M. Wu: "Safety verification of deep neural networks", in International Conference on Computer Aided Verification. Springer, 2017, pp. 3-29.
- [i.59] O. Bastani, Y. Ioannou, L. Lampropoulos, D. Vytiniotis, A. Nori and A. Criminisi: "Measuring neural net robustness with constraints", in Advances in neural information processing systems, 2016, pp. 2613- 2621.
- [i.60] V. Tjeng, K. Y. Xiao, and R. Tedrake: "Evaluating robustness of neural networks with mixed integer programming", in International Conference on Learning Representations (ICLR), 2019.
- [i.61] S. Wang, K. Pei, J. Whitehouse, J. Yang and S. Jana: "Efficient formal safety analysis of neural networks", in Advances in Neural Information Processing Systems, 2018, pp. 6367-6377.
- [i.62] Timon Gehr, Matthew Mirman, Dana Drachler-Cohen, Petar Tsankov, Swarat Chaudhuri, Martin Vechev: "AI2: Safety and Robustness Certification of Neural Networks with Abstract Interpretation", IEEE S&P 2018.
- [i.63] G. Singh, T. Gehr, M. Püschel, and M. Vechev: "An Abstract Domain for Certifying Neural Networks", Proceedings of the ACM on Programming Languages, vol. 3, no. POPL, p. 41, 2019.
- [i.64] Yizhak Yisrael Elboher, Justin Gottschlich, Guy Katz: "An Abstraction-Based Framework for Neural Network Verification", International Conference on Computer Aided Verification 2020.
- [i.65] Guy Katz, Derek A. Huang Duligur Ibeling, Kyle Julian, Christopher Lazarus, Rachel Lim, Parth Shah, Shantanu Thakoor, Haoze Wu, Aleksandar Zeljić, David L. Dill, Mykel J. Kochenderfer, Clark Barrett: "The Marabou Framework for Verification and Analysis of Deep Neural Networks", International Conference on Computer Aided Verification 2019.
- [i.66] ERAN: Neural Network Verification Framework.
ETSI GR SAI 005 V1.1.1 (2021-03)
NOTE: Available at <https://github.com/eth-sri/eran-standards/sist/b68ecdff-bbec-48c0-a490-35a21b959449/etsi-gr-sai-005-v1-1-1-2021-03>
- [i.67] Lily Weng, Pin-Yu Chen, Lam Nguyen, Mark Squillante, Akhilan Boopathy, Ivan Oseledets, Luca Daniel: "PROVEN: Verifying Robustness of Neural Networks with a Probabilistic Approach", ICML 2019.
- [i.68] Linyi Li, Xiangyu Qi, Tao Xie, and Bo Li: "SoK: Certified Robustness for Deep Neural Networks", ArXiv, 2020.
- NOTE: Available at <https://arxiv.org/pdf/2009.04131.pdf>.
- [i.69] Shixin Tian, Guolei Yang, Ying Cai: "Detecting Adversarial Examples Through Image Transformation", AAAI 2018.
- [i.70] Weilin Xu, David Evans, Yanjun Qi: "Feature squeezing: Detecting adversarial examples in Deep Neural Networks", Network and Distributed Systems Security Symposium 2018.
- [i.71] Bo Huang, Yi Wang and Wei Wang: "Model-Agnostic Adversarial Detection by Random Perturbations", International Joint Conference on Artificial Intelligence (IJCAI-19).
- [i.72] Xin Li, Fuxin Li: "Adversarial Examples Detection in Deep Networks with Convolutional Filter Statistics", ICCV 2017.
- [i.73] Kevin Roth, Yannic Kilcher, Thomas Hofmann: "The Odds are Odd: A Statistical Test for Detecting Adversarial Example", ICML 2019.
- [i.74] Dongyu Meng, Hao Chen: "MagNet: A Two Pronged Defense against adversarial examples", ACM Conference on Computer and Communications Security (CCS) 2017.

- [i.75] Nicholas Carlini, David Wagner: "MagNet and "Efficient Defenses Against Adversarial Attacks" are Not Robust to Adversarial Examples", arXiv 2017.
- NOTE: Available at <https://arxiv.org/abs/1711.08478>.
- [i.76] Faiq Khalid, Hassan Ali, Hammad Tariq, Muhammad Abdullah Hanif, Semeen Rehman, Rehan Ahmed, Muhammad Shafique: "QuSecNets: Quantization-based Defense Mechanism for Securing Deep Neural Network against Adversarial Attacks", IEEE 25th International Symposium on On-Line Testing and Robust System Design (IOLTS) 2019.
- [i.77] [Sanjay Kariyappa](#), [Moinuddin K. Qureshi](#): "Improving adversarial robustness of ensembles with diversity training", ArXiv 2019.
- NOTE: Available at <https://arxiv.org/abs/1901.09981>.
- [i.78] Thilo Strauss, Markus Hanselmann, Andrej Junginger, Holger Ulmer: "Ensemble Methods as a Defense to Adversarial Perturbations Against Deep Neural Networks", Canadian Conference on Artificial Intelligence 2020.
- NOTE: Available at http://doi-org-443.webvpn.fjmu.edu.cn/10.1007/978-3-030-47358-7_1.
- [i.79] Jialong Zhang, Zhongshu Gu, Jiyong Jang, Hui Wu, Marc Ph. Stoecklin, Heqing Huang, Ian Molloy: "Protecting Intellectual Property of Deep Neural Networks with Watermarking", ASIACCS'18.
- NOTE: Available at <https://dl.acm.org/doi/pdf/10.1145/3196494.3196550>.
- [i.80] Wang, Tianhao, and Florian Kerschbaum: "Attacks on digital watermarks for deep neural networks", in ICASSP 2019-2019 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP), pp. 2622-2626. IEEE, 2019.
- [i.81] Takemura, Tatsuya, Naoto Yandai, and Toru Fujiwara: "Model Extraction Attacks against Recurrent Neural Networks", arXiv preprint arXiv:2002.00123 (2020).
- NOTE: Available at <https://arxiv.org/abs/2002.00123>.
- [i.82] Manish Kesarwani, Bhaskar Mukhoty, Vijay Arya, Sameep Mehta: "Model Extraction Warning in MLaaS Paradigm", ACSAC 2018.
- [i.83] Huadi Zheng, Qingqing Ye, Haibo Hu, Chengfang Fang, Jie Shi: "BDPL: A Boundary Differentially Private Layer Against Machine Learning Model Extraction Attacks", ESORICS 2019.
- [i.84] Mika Juuti, Sebastian Szyller, Samuel Marchal, N. Asokan: "PRADA: Protecting Against DNN Model Stealing Attacks", 2019 IEEE European Symposium on Security and Privacy (EuroS&P).
- [i.85] Florian Tramèr, Fan Zhang, Ari Juels, Michael K. Reiter, Thomas Ristenpart: "Stealing Machine Learning Models via Prediction APIs", Usenix Security 2016.
- [i.86] Liu, Yuntao, Dana Dachman-Soled, and Ankur Srivastava: "Mitigating reverse engineering attacks on deep neural networks", in 2019 IEEE Computer Society Annual Symposium on VLSI (ISVLSI), pp. 657-662.
- [i.87] Lukas, Nils, Yuxuan Zhang, and Florian Kerschbaum: "Deep Neural Network Fingerprinting by Conferrable Adversarial Examples", arXiv preprint arXiv:1912.00888v3, (2020).
- NOTE: Available at <https://arxiv.org/abs/1912.00888v3>.
- [i.88] Cao, Xiaoyu, Jinyuan Jia, and Neil Zhenqiang Gong: "IPGuard: Protecting the Intellectual Property of Deep Neural Networks via Fingerprinting the Classification Boundary", ACM ASIA Conference on Computer and Communications Security (ASIACCS), 2021.
- [i.89] Nicolas Papernot, Martin Abadi, Ulfar Erlingsson, Ian Goodfellow, Kunal Talwar: "Semi-supervised knowledge transfer for deep learning from private training data", ICLR 2017.

- [i.90] Max Friedrich, Arne Köhn, Gregor Wiedemann, Chris Biemann: "Adversarial learning of privacy preserving test representations for de-identification of medical records", Proceeding of [the 57th Annual Meeting of the Association for Computational Linguistics](#), 2019.

NOTE: Available at <https://doi.org/10.18653/v1/P19-1584>.

- [i.91] Taihong Xiao, Yi-Hsuan Tsai, Kihyuk Sohn, Manmohan Chandraker, Ming-Hsuan Yang: "Adversarial Learning of Privacy-Preserving and Task-Oriented Representations", AAAI 2020.
- [i.92] Cynthia Dwork, Aaron Roth: "The algorithmic foundations of differential privacy", Foundations and Trends in Theoretical Computer Science, 2014.
- [i.93] Martín Abadi, Andy Chu, Ian Goodfellow, H. Brendan McMahan, Ilya Mironov, Kunal Talwar, Li Zhang: "Deep learning with differential privacy", Proceedings of the 2016 ACM CCS, 2016.
- [i.94] Milad Nasr, Reza Shokri, and Amir Houmansadr: "Machine learning with membership privacy using adversarial regularization", Proceedings of the 2018 ACM CCS, 2018.
- [i.95] Jinyuan Jia, Ahmed Salem, Michael Backes, Yang Zhang, Neil Zhenqiang Gong: "MemGuard: Defending against Black-Box Membership Inference Attacks via Adversarial Examples", ACM CCS 2019.
- [i.96] Ziqi Yang, Bin Shao, Bohan Xuan, Ee-Chien Chang, Fan Zhang: "Defending Model Inversion and Membership Inference Attacks via Prediction Purification", ArXiv, arXiv:2005.03915v2, 2020.

NOTE: Available at <https://arxiv.org/abs/2005.03915v2>.

3 Definition of terms, symbols and abbreviations

3.1 Terms

For the purposes of the present document, the following terms apply:

adversarial examples: carefully crafted samples which mislead a model to give an incorrect prediction

conferrable adversarial examples: subclass of transferable adversarial examples that exclusively transfer with a target label from a source model to its surrogates

distributional shift: distribution of input data changes over time

inference attack: attacks launched from deployment stage

model-agnostic mitigation: mitigations which do not modify the addressed machine learning model

model enhancement mitigation: mitigations which modify the addressed machine learning model

training attack: attacks launched from development stage

transferable adversarial examples: adversarial examples which are crafted for one model but also fool a different model with a high probability

3.2 Symbols

Void.

3.3 Abbreviations

For the purposes of the present document, the following abbreviations apply:

AE	Adversarial Example
AI	Artificial Intelligence
API	Application Interface
BDP	Boundary Differential Privacy
BIM	Basic Iterative Method
CNN	Convolutional Neural Network
CW	Carlini & Wagner (attacks)
DNN	Deep Neural Network
DP-SGD	Differential-Privacy Stochastic Gradient Descent
FGSM	Fast Gradient Sign Method
GNN	Graph Neural Network
ILSVRC	ImageNet Large Scale Visual Recognition Challenge
IP	Intellectual Property
JPEG	Joint Photographic Experts Group
JSMA	Jacobian-based Saliency Map
KNHT	Keyed Non-parametric Hypothesis Tests
L-BFGS	Limited-Memory Broyden–Fletcher–Goldfarb–Shanno (algorithm)
ML	Machine Learning
MNIST	Modified National Institute of Standards and Technology
MNTD	Meta Neural Trojan Detection
PATE	Private Aggregation of Teacher Ensemble
PCA	Principal Component Analysis
PGD	Project Gradient Descent
PRADA	Protecting Against DNN Model Stealing Attacks
ReLU	Rectified Linear Unit
RNN	Recurrent Neural Network
RONI	Reject On Negative Impact
SAI	Securing Artificial Intelligence
SAT	Satisfiability
SGD	Stochastic Gradient Descent
SMT	Satisfiability Modulo Theories
STRIP	STRong Intentional Perturbation
TRIM	Trimmed-based algorithm
ULP	Universal Litmus Pattern

4 Overview

4.1 Machine learning models workflow

Artificial intelligence has been driven by the rapid progress in deep learning and the wide applications of deep learning, such as image classification, object detection, speech recognition and language translation. Therefore, the present document focuses on deep learning and explores existing mitigations countermeasuring attacks on deep learning.

A machine learning model workflow is represented in Figure 1. The model life-cycle includes both development and deployment stages. The *training dataset* is the subset of domain data samples used to train the model, and it can be obtained from one or multiple data sources, represented in Figure 1 as *data supply chain*. A pretrained model can be used as input to create the target model. At development stage, via the training dataset, the model is trained. The trained model is then tested. Pursuant to ETSI GR SAI 004 [i.1], the testing step will include functional test and adversarial test. At deployment stage, the trained and tested model is deployed, i.e. becomes the *model in operation*. Given *inference input*, the *model in operation* delivers an *output*. In Figure 1, the dotted lines from the *model in operation* back to the *model under development* capture the model *updates* in online learning scenarios [i.2]. *Updates* can be pairs of inference input and user feedback, served as new training data to refine the model. *Updates* can also be locally-computed model parameter refinements. These multiple dotted lines between the *model under development* and the *model in operation* capture the federated learning scenarios, where a global model is distributed among several entities and entities provide model updates to refine the global model.